## CLAIMS

1    1.    A method for speech synthesis, comprising:

2        providing a segment inventory comprising, for a
3    plurality of speech segments, respective sequences of
4    feature vectors, by estimating spectral envelopes of
5    input speech signals corresponding to the speech segments
6    in a succession of time intervals during each of the
7    speech segments, and integrating the spectral envelopes
8    over a plurality of window functions in a frequency
9    domain so as to determine vector elements of the feature
10   vectors;

11       receiving phonetic and prosodic information
12   indicative of an output speech signal to be generated;

13       selecting the sequences of feature vectors from the
14   inventory responsive to the phonetic and prosodic
15   information;

16       processing the selected sequences of feature vectors
17   so as to generate a concatenated output series of feature
18   vectors;

19       computing a series of complex line spectra of the
20   output signal from the series of the feature vectors; and

21       transforming the complex line spectra to a time
22   domain speech signal for output.

1    2.    A method according to claim 1, wherein providing the
2    segment inventory comprises providing segment information
3    comprising respective phonetic identifiers of the
4    segments, and wherein selecting the sequences of feature
5    vectors comprises finding the segments whose phonetic
6    identifiers are close to the received phonetic
7    information.

1  3.   A method according to claim 2, wherein the segments
2  comprise lefemes, and wherein the phonetic identifiers
3  comprise lefeme labels.

1  4.   A method according to claim 2, wherein the segment
2  information further comprises one or more prosodic
3  parameters with respect to each of the segments, and
4  wherein selecting the sequences of feature vectors
5  comprises finding the segments whose one or more prosodic
6  parameters are close to the received prosodic
7  information.

1  5.   A method according to claim 4, wherein the one or
2  more prosodic parameters are selected from a group of
3  parameters consisting of a duration, an energy level and
4  a pitch of each of the segments.

1  6.   A method according to claim 1, wherein the feature
2  vectors comprise auxiliary vector elements indicative of
3  further features of the speech segments, in addition to
4  the elements determined by integrating the spectral
5  envelopes of the input speech signals.

1  7.   A method according to claim 6, wherein the auxiliary
2  vector elements comprise voicing vector elements
3  indicative of a degree of voicing of frames of the
4  corresponding speech segments, and wherein computing the
5  complex line spectra comprises reconstructing the output
6  speech signal with the degree of voicing indicated by the
7  voicing vector elements.

1  8.   A method according to claim 7, wherein receiving the
2  prosodic information comprises receiving pitch values,
3  and wherein reconstructing the output speech signal

4   comprises adjusting a frequency spectrum of the output
5   speech signal responsive to the pitch values.

1   9.   A method according to claim 1, wherein selecting the
2   sequences of feature vectors comprises:

3      selecting candidate segments from the inventory;

4      computing a cost function for each of the candidate
5   segments responsive to the phonetic and prosodic
6   information and to the feature vectors of the candidate
7   segments; and

8      selecting the segments so as to minimize the cost
9   function.

1   10.   A method according to claim 1, wherein concatenating
2   the selected sequences of feature vectors comprises
3   adjusting the feature vectors responsive to the prosodic
4   information.

1   11.   A method according to claim 10, wherein the prosodic
2   information comprises respective durations of the
3   segments to be incorporated in the output speech signal,
4   and wherein adjusting the feature vectors comprises
5   removing one or more of the feature vectors from the
6   selected sequences so as to shorten the durations of one
7   or more of the segments.

1   12.   A method according to claim 10, wherein the prosodic
2   information comprises respective durations of the
3   segments to be incorporated in the output speech signal,
4   and wherein adjusting the feature vectors comprises
5   adding one or more further feature vectors to the
6   selected sequences so as to lengthen the durations of one
7   or more of the segments.

1   13.  A method according to claim 10, wherein the prosodic
2   information comprises respective energy levels of the
3   segments to be incorporated in the output speech signal,
4   and wherein adjusting the feature vectors comprises
5   altering one or more of the vector elements so as to
6   adjust the energy levels of one or more of the segments.

1   14.  A method according to claim 1, wherein processing
2   the selected sequences comprises adjusting the vector
3   elements so as to provide a smooth transition between the
4   segments in the time domain signal.

1   15.  A method according to claim 1, wherein the vector
2   elements comprise Mel Frequency Cepstral Coefficients of
3   the speech segments, determined based on the integrated
4   spectral envelopes.

1   16.  A method for speech synthesis, comprising:
2       receiving an input speech signal containing a set of
3   speech segments;
4       estimating spectral envelopes of the input speech
5   signal in a succession of time intervals during each of
6   the speech segments;
7       integrating the spectral envelopes over a plurality
8   of window functions in a frequency domain so as to
9   determine elements of feature vectors corresponding to
10  the speech segments; and
11      reconstructing an output speech signal by
12  concatenating the feature vectors corresponding to a
13  sequence of the speech segments.

1   17.  A method according to claim 16, wherein receiving
2   the input speech signal comprises dividing the input
3   speech signal into the segments and determining segment

4 information comprising respective phonetic identifiers of
5 the segments, and wherein reconstructing the output
6 speech signal comprises selecting the segments whose
7 feature vectors are to be concatenated responsive to the
8 segment information determined with respect to the
9 segments.

1 18. A method according to claim 17, wherein dividing the
2 input speech signal into the segments comprises dividing
3 the signal into lefemes, and wherein the phonetic
4 identifiers comprise lefeme labels.

1 19. A method according to claim 17, wherein determining
2 the segment information further comprises finding
3 respective segment parameters including one or more of a
4 duration, an energy level and a pitch of each of the
5 segments, responsive to which parameters the segments are
6 selected for use in reconstructing the output speech
7 signal.

1 20. A method according to claim 19, wherein
2 reconstructing the output speech signal comprises
3 modifying the feature vectors of the selected segments so
4 as to adjust the segment parameters of the segments in
5 the output speech signal.

1 21. A method according to claim 16, and comprising
2 determining respective degrees of voicing of the speech
3 segments, and incorporating the degrees of voicing as
4 elements of the feature vectors for use in reconstructing
5 the output speech signal.

1 22. A method according to claim 16, wherein
2 concatenating the feature vectors comprises concatenating
3 the vectors to form a series in a frequency domain, and

4 wherein reconstructing the output speech signal comprises
5 computing a series of complex line spectra of the output
6 signal from the series of feature vectors, and
7 transforming the complex line spectra to a time domain
8 signal.

1 23. A method according to claim 16, wherein the window
2 functions are non-zero only within different, respective
3 spectral windows and have variable values over their
4 respective windows, and wherein integrating the spectral
5 envelopes comprises calculating products of the spectral
6 envelopes with the window functions, and calculating
7 integrals of the products over the respective windows of
8 the window functions.

1 24. A method according claim 23, and comprising applying
2 a mathematical transformation to the integrals in order
3 to determine the elements of the feature vectors.

1 25. A method according to claim 24, wherein the
2 frequency domain comprises a Mel frequency domain, and
3 wherein applying the mathematical transformation
4 comprises applying log and discrete cosine transform
5 operations in order to determine Mel Frequency Cepstral
6 Coefficients to be used as the elements of the feature
7 vectors.

1 26. A device for speech synthesis, comprising:
2 a memory, arranged to hold a segment inventory
3 comprising, for a plurality of speech segments,
4 respective sequences of feature vectors having vector
5 elements determined by estimating spectral envelopes of
6 input speech signals corresponding to the speech segments
7 in a succession of time intervals during each of the

8  speech segments, and integrating the spectral envelopes
9  over a plurality of window functions in a frequency
10 domain; and

11      a speech processor, arranged to receive phonetic and
12 prosodic information indicative of an output speech
13 signal to be generated, to select the sequences of
14 feature vectors from the inventory responsive to the
15 phonetic and prosodic information, to process the
16 selected sequences of feature vectors so as to generate a
17 concatenated output series of feature vectors, and to
18 compute a series of complex line spectra of the output
19 signal from the series of the feature vectors and
20 transform the complex line spectra to a time domain
21 speech signal for output.

1  27.  A device according to claim 26, wherein the segment
2  inventory comprises segment information comprising
3  respective phonetic identifiers of the segments, and
4  wherein the processor is arranged to select the sequences
5  of feature vectors by finding the segments in the
6  inventory whose phonetic identifiers are close to the
7  received phonetic information.

1  28.  A device according to claim 27, wherein the segments
2  comprise lefemes, and wherein the phonetic identifiers
3  comprise lefeme labels.

1  29.  A device according to claim 27, wherein the segment
2  information further comprises one or more prosodic
3  parameters with respect to each of the segments, and
4  wherein the processor is arranged to select the sequences
5  of feature vectors by finding the segments whose one or
6  more prosodic parameters are close to the received
7  prosodic information.

1 30. A device according to claim 29, wherein the one or
2 more prosodic parameters are selected from a group of
3 parameters consisting of a duration, an energy level and
4 a pitch of each of the segments.

1 31. A device according to claim 26, wherein the feature
2 vectors comprise auxiliary vector elements indicative of
3 further features of the speech segments, in addition to
4 the elements determined by integrating the spectral
5 envelopes of the input speech signals.

1 32. A device according to claim 31, wherein the
2 auxiliary vector elements comprise voicing vector
3 elements indicative of a degree of voicing of frames of
4 the corresponding speech segments, and wherein the
5 processor is arranged to reconstruct the output speech
6 signal with the degree of voicing indicated by the
7 voicing vector elements.

1 33. A device according to claim 32, wherein the prosodic
2 information comprises pitch values, and wherein the
3 processor is arranged to adjust a frequency spectrum of
4 the output speech signal responsive to the pitch values.

1 34. A device according to claim 26, wherein the
2 processor is arranged to select the sequences of feature
3 ·vectors by selecting candidate segments from the
4 inventory, computing a cost function for each of the
5 candidate segments responsive to the phonetic and
6 prosodic information and to the feature vectors of the
7 candidate segments, and selecting the segments so as to
8 minimize the cost function.

1 35. A device according to claim 26, wherein the
2 processor is arranged to adjust the feature vectors in

3  the combined output series responsive to the prosodic
4  information.

1  36. A device according to claim 35, wherein the prosodic
2  information comprises respective durations of the
3  segments to be incorporated in the output speech signal,
4  and wherein the processor is arranged to adjust the
5  feature vectors by removing one or more of the feature
6  vectors from the selected sequences so as to shorten the
7  durations of one or more of the segments.

1  37. A device according to claim 35, wherein the prosodic
2  information comprises respective durations of the
3  segments to be incorporated in the output speech signal,
4  and wherein the processor is arranged to adjust the
5  feature vectors by adding one or more further feature
6  vectors to the selected sequences so as to lengthen the
7  durations of one or more of the segments.

1  38. A device according to claim 35, wherein the prosodic
2  information comprises respective energy levels of the
3  segments to be incorporated in the output speech signal,
4  and wherein the processor is arranged to adjust the
5  energy levels of one or more of the segments by altering
6  one or more of the vector elements.

1  39. A device according to claim 26, wherein the
2  processor is arranged to adjust the vector elements so as
3  to provide a smooth transition between the segments in
4  the time domain signal.

1  40. A device according to claim 26, wherein the vector
2  elements comprise Mel Frequency Cepstral Coefficients of
3  the speech segments, determined based on the integrated
4  spectral envelopes.

1 41. A device for speech synthesis, comprising:

2 a memory, arranged to hold a segment inventory
3 determined by processing an input speech signal
4 containing a set of speech segments so as to estimate
5 spectral envelopes of the input speech signal in a
6 succession of time intervals during each of the speech
7 segments, and integrating the spectral envelopes over a
8 plurality of window functions in a frequency domain so as
9 to determine elements of feature vectors corresponding to
10 the speech segments; and

11 a speech processor, arranged to reconstruct an
12 output speech signal by concatenating the feature vectors
13 corresponding to a sequence of the speech segments.

1 42. A device according to claim 41, wherein the input
2 speech signal is processed by dividing the input speech
3 signal into the segments and determining segment
4 information comprising respective phonetic identifiers of
5 the segments, and wherein the processor is arranged to
6 reconstruct the output speech signal by selecting the
7 segments whose feature vectors are to be concatenated
8 responsive to the segment information determined with
9 respect to the segments.

1 43. A device according to claim 42, wherein the input
2 speech signal is divided into lefemes, and the phonetic
3 identifiers comprise lefeme labels.

1 44. A device according to claim 42, wherein the segment
2 information further comprises respective segment
3 parameters including one or more of a duration, an energy
4 level and a pitch of each of the segments, responsive to
5 which parameters the segments are selected by the

6   processor for use in reconstructing the output speech
7   signal.

1   45. A device according to claim 44, wherein the
2   processor is arranged to modify the feature vectors of
3   the selected segments so as to adjust the segment
4   parameters of the segments in the output speech signal.

1   46. A device according to claim 41, wherein the feature
2   vectors comprise respective degrees of voicing of the
3   speech segments, for use by the processor in
4   reconstructing the output speech signal.

1   47. A device according to claim 41, wherein the
2   processor is arranged to concatenate the feature vectors
3   to form a series in a frequency domain, and to
4   reconstruct the output speech signal by computing a
5   series of complex line spectra of the output signal from
6   the series of feature vectors, and transforming the
7   complex line spectra to a time domain signal.

1   48. A device according to claim 14, wherein the window
2   functions are non-zero only within different, respective
3   spectral windows and have variable values over their
4   respective windows, and wherein the feature vector
5   elements are determined by calculating products of the
6   spectral envelopes with the window functions, and
7   calculating integrals of the products over the respective
8   windows of the window functions.

1   49. A device according claim 48, wherein a mathematical
2   transformation is applied to the integrals in order to
3   determine the elements of the feature vectors.

1   50. A device according to claim 48, wherein the
2   frequency domain comprises a Mel frequency domain, and

3   wherein the mathematical transformation comprises log and
4   discrete cosine transform operations, which are applied
5   so as to determine Mel Frequency Cepstral Coefficients to
6   be used as the elements of the feature vectors.

1   51. A computer software product, comprising a
2   computer-readable medium in which program instructions
3   are stored, which instructions, when read by a computer,
4   cause the computer to access a segment inventory
5   comprising, for a plurality of speech segments,
6   respective sequences of feature vectors having vector
7   elements determined by estimating spectral envelopes of
8   input speech signals corresponding to the speech segments
9   in a succession of time intervals during each of the
10   speech segments, and integrating the spectral envelopes
11   over a plurality of window functions in a frequency
12   domain, and in response to phonetic and prosodic
13   information indicative of an output speech signal to be
14   generated, cause the computer to select the sequences of
15   feature vectors from the inventory responsive to the
16   phonetic and prosodic information, to process the
17   selected sequences of feature vectors so as to generate a
18   concatenated output series of feature vectors, and to
19   compute a series of complex line spectra of the output
20   signal from the series of the feature vectors and
21   transform the complex line spectra to a time domain
22   speech signal for output.

1   52. A product according to claim 51, wherein the segment
2   inventory comprises segment information comprising
3   respective phonetic identifiers of the segments, and
4   wherein the instructions cause the computer to select the
5   sequences of feature vectors by finding the segments in

6  the inventory whose phonetic identifiers are close to the
7  received phonetic information.

1  53. A product according to claim 52, wherein the
2  segments comprise lefemes, and wherein the phonetic
3  identifiers comprise lefeme labels.

1  54. A product according to claim 52, wherein the segment
2  information further comprises one or more prosodic
3  parameters with respect to each of the segments, and
4  wherein the instructions cause the computer to select the
5  sequences of feature vectors by finding the segments
6  whose one or more prosodic parameters are close to the
7  received prosodic information.

1  55. A product according to claim 54, wherein the one or
2  more prosodic parameters are selected from a group of
3  parameters consisting of a duration, an energy level and
4  a pitch of each of the segments.

1  56. A product according to claim 54, wherein the feature
2  vectors comprise auxiliary vector elements indicative of
3  further features of the speech segments, in addition to
4  the elements determined by integrating the spectral
5  envelopes of the input speech signals.

1  57. A product according to claim 56, wherein the
2  auxiliary vector elements comprise voicing vector
3  elements indicative of a degree of voicing of frames of
4  the corresponding speech segments, and wherein the
5  instructions cause the computer to reconstruct the output
6  speech signal with the degree of voicing indicated by the
7  voicing vector elements.

1  58. A product according to claim 57, wherein the
2  prosodic information comprises pitch values, and wherein

3  the instructions cause the computer to adjust a frequency
4  spectrum of the output speech signal responsive to the
5  pitch values.

1  59.  A  product  according  to  claim  51,  wherein  the
2  instructions cause the computer to select the sequences
3  of feature vectors by selecting candidate segments from
4  the inventory, computing a cost function for each of the
5  candidate  segments  responsive  to  the  phonetic  and
6  prosodic information and to the feature vectors of the
7  candidate segments, and selecting the segments so as to
8  minimize the cost function.

1  60.  A  product  according  to  claim  51,  wherein  the
2  instructions cause the computer to adjust the feature
3  vectors in the combined output series responsive to the
4  prosodic information.

1  61.  A  product  according  to  claim  60,  wherein  the
2  prosodic  information  comprises  respective  durations  of
3  the  segments  to  be  incorporated  in  the  output  speech
4  signal, and wherein the instructions cause the computer
5  to adjust the feature vectors by removing one or more of
6  the feature vectors from the selected sequences so as to
7  shorten the durations of one or more of the segments.

1  62.  A  product  according  to  claim  60,  wherein  the
2  prosodic  information  comprises  respective  durations  of
3  the  segments  to  be  incorporated  in  the  output  speech
4  signal, and wherein the instructions cause the computer
5  to  adjust  the  feature  vectors  by  adding  one  or  more
6  further feature vectors to the selected sequences so as
7  to lengthen the durations of one or more of the segments.

1 63. A product according to claim 60, wherein the
2 prosodic information comprises respective energy levels
3 of the segments to be incorporated in the output speech
4 signal, and wherein the instructions cause the computer
5 to adjust the energy levels of one or more of the
6 segments by altering one or more of the vector elements.

1 64. A product according to claim 51, wherein the
2 instructions cause the computer to adjust the vector
3 elements so as to provide a smooth transition between the
4 segments in the time domain signal.

1 65. A product according to claim 51, wherein the vector
2 elements comprise Mel Frequency Cepstral Coefficients of
3 the speech segments, determined based on the integrated
4 spectral envelopes.

1 66. A computer software product, comprising a
2 computer-readable medium in which a segment inventory is
3 stored, the inventory having been determined by
4 processing an input speech signal containing a set of
5 speech segments so as to estimate spectral envelopes of
6 the input speech signal in a succession of time intervals
7 during each of the speech segments, and integrating the
8 spectral envelopes over a plurality of window functions
9 in a frequency domain so as to determine elements of
10 feature vectors corresponding to the speech segments, so
11 that a speech processor can reconstruct an output speech
12 signal by concatenating the feature vectors corresponding
13 to a sequence of the speech segments.

1 67. A product according to claim 66, wherein the input
2 speech signal is processed by dividing the input speech
3 signal into the segments and determining segment

4  information comprising respective phonetic identifiers of
5  the segments, and wherein to reconstruct the output
6  speech signal, the processor selects the segments whose
7  feature vectors are to be concatenated responsive to the
8  segment information determined with respect to the
9  segments.

1  68. A product according to claim 66, wherein the input
2  speech signal is divided into lefemes, and the phonetic
3  identifiers comprise lefeme labels.

1  69. A product according to claim 66, wherein the segment
2  information further comprises respective segment
3  parameters including one or more of a duration, an energy
4  level and a pitch of each of the segments, responsive to
5  which parameters the segments are selected by the
6  computer for use in reconstructing the output speech
7  signal.

1  70. A product according to claim 69, wherein to
2  reconstruct the output speech signal, the instructions
3  cause the computer to modify the feature vectors of the
4  selected segments so as to adjust the durations and
5  energy levels of the segments in the output speech
6  signal.

1  71. A product according to claim 66, wherein the feature
2  vectors comprise respective degrees of voicing of the
3  speech segments, for use by the computer in
4  reconstructing the output speech signal.

1  72. A product according to claim 66, wherein to
2  reconstruct the output speech signal, the instructions
3  cause the computer to concatenate the feature vectors to
4  form a series in a frequency domain, to compute as series

5  of complex line spectra of the output signal from the
6  series of feature vectors, and to transform the complex
7  line spectra to a time domain signal.

1  73.  A product according to claim 66, wherein the window
2  functions are non-zero only within different, respective
3  spectral windows and have variable values over their
4  respective windows, and wherein the feature vector
5  elements are determined by calculating products of the
6  spectral envelopes with the window functions, and
7  calculating integrals of the products over the respective
8  windows of the window functions.

1  74.  A product according claim 73, wherein a mathematical
2  transformation is applied to the integrals in order to
3  determine the elements of the feature vectors.

1  75.  A product according to claim 74, wherein the
2  frequency domain comprises a Mel frequency domain, and
3  wherein the mathematical transformation comprises log and
4  discrete cosine transform operations, which are applied
5  so as to determine Mel Frequency Cepstral Coefficients to
6  be used as the elements of the feature vectors.